

Chapter 8: Bootstrapping

Typically in statistics, we use **theory** to derive the sampling distribution of a statistic. From the sampling distribution, we can obtain the variance, construct confidence intervals, perform hypothesis tests, and more.

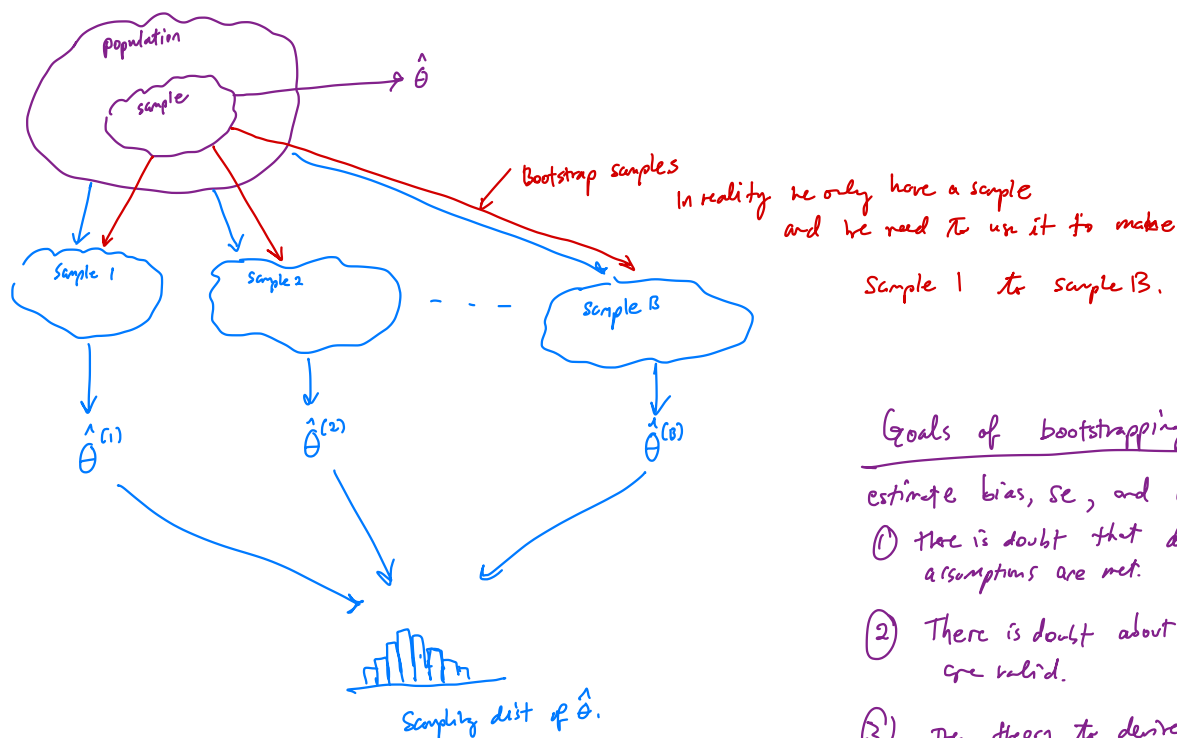
Challenge:

What if the sampling distribution is impossible to obtain or asymptotic theory doesn't hold?

Basic idea of bootstrapping: "Pull yourself up by your bootstraps"

- Use the data to estimate the sampling distribution of the statistic.
- Estimate the sampling distribution by [creating a large number of datasets that we might have seen] and compute the statistic on each of these data sets.

Eg:



Goals of bootstrapping

- estimate bias, SE, and CI's when
- ① there is doubt that distributional assumptions are met.
 - ② There is doubt about asymptotic results are valid.
 - ③ the theory to derive the den of the test statistic is too hard.

1 Nonparametric Bootstrap

Let $X_1, \dots, X_n \sim F$ with pdf $f(x)$. Recall, the cdf is defined as

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt.$$

Definition 1.1 The empirical cdf is a function which estimates the cdf using observed data,

$$\hat{F}(x) = F_n(x) = \text{proportion of sample points that fall in } (-\infty, x].$$

~ "depends on the data"

In practice, this leads to the following function. Let $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ be the order statistics of the sample. Then,

= sample in order

$$F_n(x) = \begin{cases} 0 & x < X_{(1)} \\ \frac{i}{n} & X_{(i)} \leq x < X_{(i+1)}; \quad i = 1, \dots, n-1 \\ 1 & x \geq X_{(n)} \end{cases}$$

$F_n(x)$ is an estimator of cdf $F(x)$ and as $n \rightarrow \infty$, $F_n \rightarrow F$

Theoretical: Sample $X \sim F$ use X_1, \dots, X_n to compute F_n .

Bootstrap: Sample $X^* \sim F_n$, use X_1^*, \dots, X_n^* to compute F_n^*

Example 1.1 Let $x = 2, 2, 1, 1, 5, 4, 4, 3, 1, 2$ be an observed sample. Find $F_n(x)$.

Sorted $x = 1, 1, 1, 2, 2, 2, 3, 4, 4, 5$ $n=10$.

$$F_n(x) = \begin{cases} 0 & x < 1 \\ \frac{3}{10} & 1 \leq x < 2 \\ \frac{6}{10} & 2 \leq x < 3 \\ \frac{7}{10} & 3 \leq x < 4 \\ \frac{9}{10} & 4 \leq x < 5 \\ 1 & x \geq 5 \end{cases}$$

There is an easy way to sample from F_n without calculating it.

The idea behind the bootstrap is to sample many data sets from $F_n(x)$, which can be achieved by resampling from the data with replacement.

```
# observed data
x <- c(2, 2, 1, 1, 5, 4, 4, 3, 1, 2)
```

```
# create 10 bootstrap samples
```

```
x_star <- matrix(NA, nrow = length(x), ncol = 10)
for(i in 1:10) {
  x_star[, i] <- sample(x, length(x), replace = TRUE)
}
x_star
```

key part of the bootstrap!

sample of length n from $F_n(x)$.

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    5    2    4    4    1    2    2    1    5    1
## [2,]    4    5    1    1    1    2    1    1    4    2
## [3,]    4    2    5    1    2    2    1    4    4    3
## [4,]    4    5    1    3    2    4    4    4    3    1
## [5,]    4    1    2    1    1    1    5    2    1    1
## [6,]    4    2    2    2    4    4    3    2    1    2
## [7,]    1    5    4    4    1    2    1    2    1    4
## [8,]    3    1    1    1    4    1    4    1    4    2
## [9,]    1    4    4    2    2    1    4    3    2    1
## [10,]   4    1    2    3    4    5    5    5    2    4
```

$x^{(1)}$

$x^{(10)}$

```
# compare mean of the sample to the means of the bootstrap samples
mean(x)
```

```
## [1] 2.5  $\bar{x}$   
           $(\hat{\theta})$ 
```

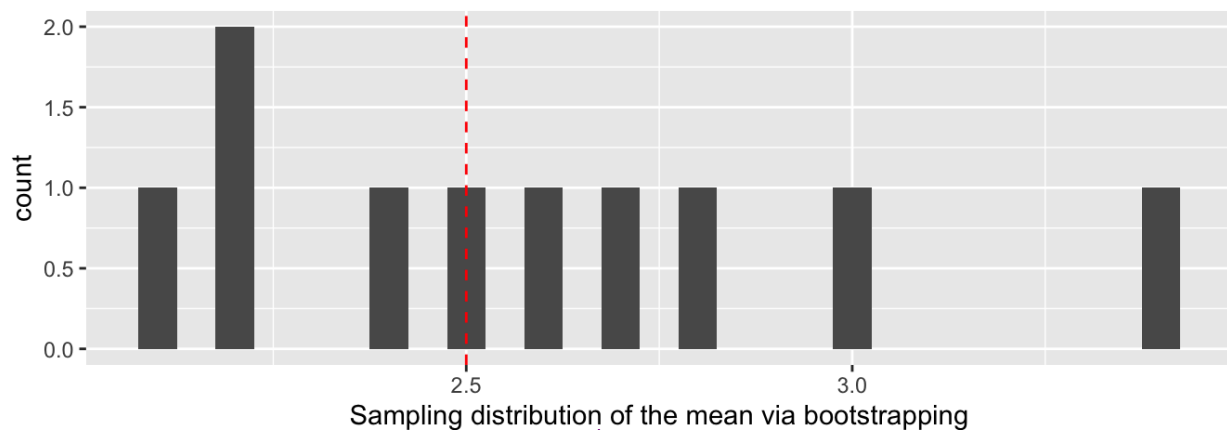
```
colMeans(x_star)
```

```
## [1] 3.4 2.8 2.6 2.2 2.2 2.4 3.0 2.5 2.7 2.1
```

$\bar{x}^{(1)}$
 $\hat{\theta}^{(1)}$

$\bar{x}^{(10)}$
 $\hat{\theta}^{(10)}$

```
ggplot() +
  geom_histogram(aes(colMeans(x_star)), binwidth = .05) +
  geom_vline(aes(xintercept = mean(x)), lty = 2, colour = "red") +
  xlab("Sampling distribution of the mean via bootstrapping")
```



1.1 Algorithm

Goal: estimate the sampling distribution of a statistic based on observed data x_1, \dots, x_n .

Let θ be the parameter of interest and $\hat{\theta} = T(x_1, \dots, x_n)$ be an estimator of θ . Then, \uparrow
n observations.

For $b = 1, \dots, B$ \leftarrow # of bootstrap samples.

① Sample $x^{*(b)} = (x_1^{*(b)}, \dots, x_n^{*(b)})$ by sample with replacement from observed data x
(i.e. sampling from F_n)

② $\hat{\theta}^{(b)} = T(x_1^{*(b)}, \dots, x_n^{*(b)})$.
 \uparrow estimate of θ based on b^{th} bootstrap sample.

Using $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}$

- estimate the sampling distribution of the statistic $\hat{\theta}$
 \hookrightarrow make a histogram of $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}$

- estimate the standard error of $\hat{\theta}$
 \hookrightarrow compute st. deviation of $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}$

- estimate a CI
 \hookrightarrow we'll cover multiple methods.

- estimate many other things as well.

1.2 Properties of Estimators

We can use the bootstrap to estimate different properties of estimators.

1.2.1 Standard Error

Recall $se(\hat{\theta}) = \sqrt{Var(\hat{\theta})}$. We can get a **bootstrap** estimate of the standard error:

$$\hat{se}(\hat{\theta}) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}^{(b)} - \bar{\hat{\theta}})^2}$$

$$\text{where } \bar{\hat{\theta}} = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^{(b)}$$

1.2.2 Bias

Recall $bias(\hat{\theta}) = E[\hat{\theta} - \theta] = E[\hat{\theta}] - \theta$.

Example 1.2

$$E[\hat{\sigma}^2] = E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right] = \left(1 - \frac{1}{n}\right) \sigma^2$$

$$\Rightarrow bias[\hat{\sigma}^2] = E[\hat{\sigma}^2] - \sigma^2 = \left(1 - \frac{1}{n}\right) \sigma^2 - \sigma^2 = -\frac{1}{n} \sigma^2$$

$$\Rightarrow \text{we use } s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad E[s^2] = \sigma^2.$$

We can get a **bootstrap** estimate of the bias:

$$\hat{bias}(\hat{\theta}) = \underbrace{\bar{\hat{\theta}}}_{\substack{\uparrow \\ \text{computed} \\ \text{from bootstrap} \\ \text{samples}}} - \underbrace{\hat{\theta}}_{\substack{\uparrow \\ \text{based} \\ \text{on original} \\ \text{data sample}}} = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}^{(b)} - \hat{\theta}).$$

If $\hat{bias}(\hat{\theta}) > 0$, then $\hat{\theta}$ overestimate θ on average.

Overall, we seek statistics with small se and small bias.

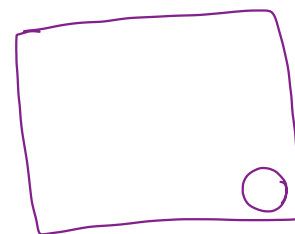
but there typically is a bias/variance tradeoff as bias \downarrow , se \uparrow

1.3 Sample Size and # Bootstrap Samples

n = sample size & B = # bootstrap samples

If n is too small, or sample isn't representative of the population,

Then bootstrap results will be poor no matter how large B is.



Guidelines for B –

$B \approx 1000$ for se & bias

$B \approx 2000$ for CI's (depends on α : small $\alpha \Rightarrow \uparrow B$)

Best approach –

Repeat bootstrap twice w/ different seeds

If estimates are very different, $\uparrow B$.

Your Turn

In this example, we explore bootstrapping in the rare case where we know the values for the entire population. If you have all the data from the population, you don't need to bootstrap (or really, inference). It is useful to learn about bootstrapping by comparing to the truth in this example.

In the package `bootstrap` is contained the average LSAT and GPA for admission to the population of 82 USA Law schools (an old data set – there are now over 200 law schools). This package also contains a random sample of size $n = 15$ from this dataset.

```
library(bootstrap)
```

```
head(law)
```

↑ random sample of size $n = 15$.

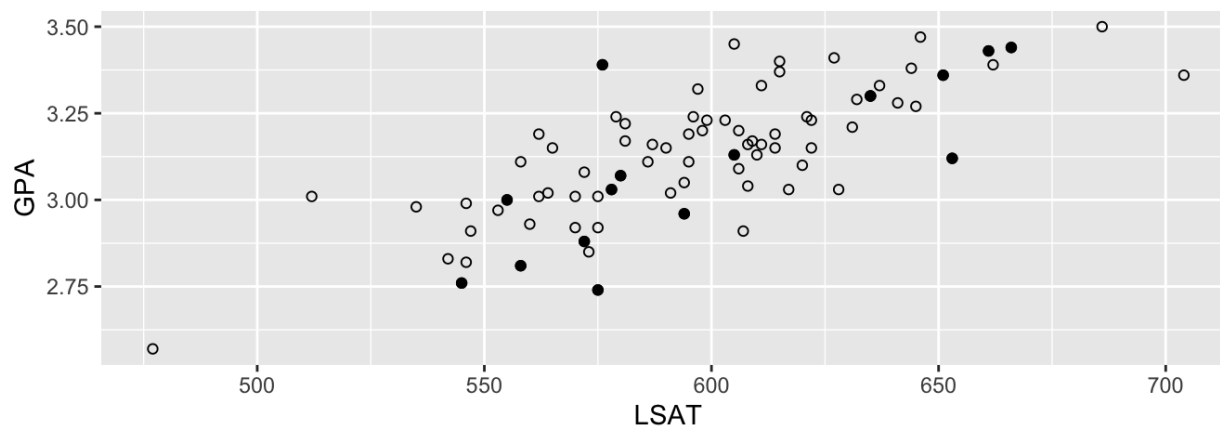
```
##   LSAT  GPA
## 1  576 3.39
## 2  635 3.30
## 3  558 2.81
## 4  578 3.03
## 5  666 3.44
## 6  580 3.07
```

```
ggplot() +
```

```
  geom_point(aes(LSAT, GPA), data = law) +
```

```
  geom_point(aes(LSAT, GPA), data = law82, pch = 1)
```

↑ Full data set (population).



$$\text{recall } \hat{\theta} = \hat{\rho} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

We will estimate the correlation $\theta = \rho(\text{LSAT}, \text{GPA})$ between these two variables and use a bootstrap to estimate the sample distribution of $\hat{\theta}$.

```
# sample correlation
cor(law$LSAT, law$GPA)
```

```
## [1] 0.7763745
```

```
# population correlation
cor(law82$LSAT, law82$GPA)
```

```
## [1] 0.7599979
```

```
# set up the bootstrap
```

```
B <- 1000
```

```
n <- nrow(law)
```

```
r <- numeric(B) # storage
```

```
for(b in B) {
  ## Your Turn: Do the bootstrap!
}
```

we know this because
we have the population.

1. Plot the sample distribution of $\hat{\theta}$. Add vertical lines for the true value θ and the sample estimate $\hat{\theta}$.
2. Estimate $se(\hat{\theta})$.
3. Estimate the bias of $\hat{\theta}$